

Exact Confidence Intervals on the Regression Coefficients in Multiple Regression Model with Nested Error Structure

Dong Joon Park¹⁾

Abstract

In regression model with nested error structure interval estimations on regression coefficients in different stages are proposed. Ordinary least square estimators and generalized least square estimators of the regression coefficients in this model are derived for between and within group model. The confidence intervals are derived by using independent distributional properties between regression coefficient estimators and quadratic forms obtained from the model.

1. Introduction

This article considers the multiple regression model where the responses are correlated. This model is appropriate to the data collected using two stage cluster designs, panel data analysis, or split-plot type models. Fuller and Battese(1973) presented transformation by which uncorrelated errors with constant variances were obtained and the transformations permitted the calculation of the Generalized Least Squares(GLS) estimators and their covariance matrices by Ordinary Least Squares(OLS) regression. Christensen(1987) analyzed two stage sampling data and showed OLS estimators were Best Linear Unbiased Estimators(BLUE) iff for each variable in the model, variable obtained by replacing each component with corresponding cluster average was also contained in the model. Weerakkoday and Johnson(1992) presented sufficient conditions under which OLS estimates of regressor parameters in nested error data structure were Uniformly Minimum Variance Unbiased(UMVU). Park and Burdick(1994) proposed the confidence intervals on the regression coefficients in simple regression model with one-fold nested error structure. Park(1996) proposed the distributional properties of variance components in two stage regression model. This article proposes the exact confidence intervals on the regression coefficients associated with primary and secondary units in multiple regression model with nested error structure.

1) Assistant Professor, Division of Mathematical Science, Pukyong National University, Nam-Gu, Daeyeon-3 Dong, Pusan 608-737

2. Multiple regression model with nested error structure

The multiple regression model with nested error structure is written as

$$\begin{aligned}
 Y_{ij} &= \beta_0 + \beta_1 X_{h1} + \cdots + \beta_{p_1} X_{hp_1} + \delta_i \\
 &\quad + \gamma_1 X_{ij1} + \cdots + \gamma_{p_2} X_{ijp_2} + \varepsilon_{ij} \\
 h &= 1, \dots, m_{p_1}; \quad i=1, \dots, l_1; \quad j=1, \dots, l_2
 \end{aligned} \tag{2.1}$$

where Y_{ij} is the j th observation in the i th cell(group), β_0 is an intercept term, $\beta_1, \dots, \beta_{p_1}$ are unknown parameters associated with primary units, X_{h1}, \dots, X_{hp_1} are fixed predictor variables in the primary unit, $\gamma_1, \dots, \gamma_{p_2}$ are unknown parameters associated with secondary units, $X_{ij1}, \dots, X_{ijp_2}$ are fixed predictor variables in the secondary unit, δ_i is a random error term in the primary unit, ε_{ij} is a random error term in the secondary unit, δ_i and ε_{ij} are jointly independent normal random variables with zero means and variances σ_δ^2 and σ_ε^2 , respectively. The index l_1 is the number of different combinations(cells) of levels among X_{ij} 's, i.e., $l_1 = m_1 \times m_2 \times \cdots \times m_{p_1}$, and l_2 is the number of repetitions within an i th cell. We consider the balanced case where l_2 's are same for all i 's. Since β 's, γ 's, X_{ij} 's, and X_{ijk} 's are fixed, and δ_i and ε_{ij} are random, model (2.1) is a mixed model.

The model (2.1) is written in matrix notation,

$$\underline{Y} = \underline{Z}X_1\underline{\beta} + X_2\underline{\gamma} + Z\underline{\delta} + \underline{\varepsilon} \tag{2.2}$$

where \underline{Y} is an $l_1 l_2 \times 1$ vector of observations, Z is an $l_1 l_2 \times l_1$ design matrix with 0's and 1's, i.e., $z = \bigoplus_{i=1}^{l_1} \mathbf{1}_{l_2}$, where $\mathbf{1}_{l_2}$ is an $l_2 \times 1$ column vector of 1's and \bigoplus is the direct sum operator, X_1 is an $l_1 \times (1 + p_1)$ matrix of known values with a column of 1's in the first column and p_1 columns of X_{ij} 's from the second column to the p_1 th column, $\underline{\beta}$ is a $(1 + p_1) \times 1$ vector of parameters associated with X_{ij} 's, X_2 is an $l_1 l_2 \times p_2$ matrix of known values with p_2 columns of X_{ijk} 's from the first column to the p_2 th column, $\underline{\gamma}$ is a $p_2 \times 1$ vector of parameters associated with X_{ijk} 's, $\underline{\delta}$ is an $l_1 \times 1$ vector of random error terms, and $\underline{\varepsilon}$ is an $l_1 l_2 \times 1$ vector of random error terms. Refer to Park(1996) for specific matrix forms in detail. By the assumptions in (2.1) the response variables are normally distributed

$$\underline{Y} \sim N(X\underline{\alpha}, \sigma_\delta^2 Z Z' + \sigma_\varepsilon^2 I_{l_1 l_2}) \tag{2.3}$$

where

$$X = (Z X_1 \quad X_2) \text{ and } \underline{\alpha} = \begin{pmatrix} \underline{\beta} \\ \underline{\gamma} \end{pmatrix}.$$

Since $Z'Z/l_2 = I_{l_1}$, premultiplying (2.2) by Z'/l_2 gives

$$\bar{Y} = X_1\beta + X_2^*\gamma + \delta + \underline{\varepsilon}^* \tag{2.4}$$

where $\bar{Y} = (Z' / l_2)Y$, $X_2^* = (Z' / l_2)X_2$, and $\underline{\varepsilon}^* = (Z' / l_2)\underline{\varepsilon}$.

The reductions in sums of squares of the model are attributable to fitting the primary and secondary fixed variables and expressed into the quadratic forms. Let $G_1 = (X^{*'} X^*)^{-1}$ and $G_2 = (\bar{X}_2' \bar{X}_2)^{-1}$ where $X^* = (X_1 \ X_2^*)$, $X_2^* = (Z' / l_2)X_2$, $\bar{X}_2 = W X_2$, and $W = I_{l_1 l_2} - Z Z' / l_2$. Define $H_1 = X^* G_1 X^{*'}$ and $H_2 = \bar{X}_2 G_2 \bar{X}_2'$. Now consider the quadratic forms $R_1 = \underline{Y}' (Z / l_2) (I_{l_1} - H_1) (Z' / l_2) \underline{Y}$ and $R_2 = \underline{Y}' W' (I_{l_1 l_2} - H_2) W \underline{Y}$. Under the distributional assumptions in (2.1) the quadratic forms $R_1 / (\sigma_\delta^2 + (\sigma_\varepsilon^2 / l_2))$ and $R_2 / \sigma_\varepsilon^2$ are chi-squared random variables with $l_1 - p_1 - p_2 - 1$ and $l_1 l_2 - l_1 - p_2$ degrees of freedom, respectively. In addition, two quadratic forms are independent (see Park(1996)).

3. Exact confidence intervals for regression coefficients

In order to construct confidence intervals on the regression coefficients the distributional properties of the estimators of β and γ are examined and they are summarized in theorems in this section. The exact confidence intervals for regression coefficients in the model are proposed using the theorems.

3.1 Confidence interval on γ_i using within group OLS estimator

Since $WZ = 0$, multiplying both sides of (2.2) on the left by W gives

$$\tilde{Y} = \bar{X}_2 \gamma + W \underline{\varepsilon} \tag{3.1}$$

where $\tilde{Y} = W \underline{Y}$. Notice that W is a symmetric and idempotent matrix. Since \tilde{Y} is a linear combination of \underline{Y} , \tilde{Y} is normally distributed

$$\tilde{Y} \sim N(\bar{X}_2 \gamma, \sigma_\varepsilon^2 W). \tag{3.2}$$

The OLS estimators for γ from the within model are expressed as p_2 elements of the vector

$$\hat{\gamma}_W = (\bar{X}_2' \bar{X}_2)^{-1} \bar{X}_2' \tilde{Y}. \tag{3.3}$$

The estimators $\hat{\gamma}_W$ are called the within group OLS estimators of γ because only the variation within each group is utilized. Since $\hat{\gamma}_W$ are a linear combination of \underline{Y} ,

$$\hat{\gamma}_W \sim N(\gamma, \sigma_\varepsilon^2 (\bar{X}_2' \bar{X}_2)^{-1}). \tag{3.4}$$

Theorem 3.1 Under the distributional assumptions in (2.1), $\hat{\gamma}_W$ and R_2 are independent.

Proof R_2 is written as $R_2 = \underline{Y}' W' (I_{l_1 l_2} - H_2) W \underline{Y} = \underline{Y}' (I_{l_1 l_2} - H_2) \underline{Y}$. Note that $(\overline{X}_2' \overline{X}_2)^{-1} \overline{X}_2' (\sigma_\epsilon^2 W) (I_{l_1 l_2} - H_2) = 0$ since $\overline{X}_2' W = \overline{X}_2'$, $W H_2 = H_2$ and $\overline{X}_2' H_2 = \overline{X}_2'$. It follows by Theorem 3 (Searle 1971, p59) that $\hat{\gamma}_{iW}$ and R_2 are independent.

Note that $(I_{l_1 l_2} - H_2)$ in R_2 is an idempotent matrix. Therefore, the within group OLS estimators of regression coefficients associated with secondary units in the model are t distributed from (3.4) and Theorem 3.1. That is, it follows from (3.4) that

$$\frac{\hat{\gamma}_{iW} - \gamma_i}{\sqrt{\overline{X}_2^{ii} \sigma_\epsilon^2}} \sim N(0, 1)$$

where $\hat{\gamma}_{iW}$ is the i th OLS estimator of regression coefficients in the secondary unit based on within group model and \overline{X}_2^{ii} is the i th diagonal element of matrix $(\overline{X}_2' \overline{X}_2)^{-1}$. It follows from Theorem 3.1 that

$$\frac{\hat{\gamma}_{iW} - \gamma_i}{\sqrt{\overline{X}_2^{ii} \sigma_\epsilon^2}} \div \sqrt{\frac{R_2}{\sigma_\epsilon^2 (l_1 l_2 - l_1 - p_2)}} \sim t(l_1 l_2 - l_1 - p_2). \tag{3.5}$$

Therefore, the proposed exact $1 - 2\alpha$ two-sided confidence interval on γ_i is

$$\hat{\gamma}_{iW} \pm t_{\alpha; l_1 l_2 - l_1 - p_2} \sqrt{\overline{X}_2^{ii} S_\epsilon^2} \tag{3.6}$$

where $S_\epsilon^2 = R_2 / (l_1 l_2 - l_1 - p_2)$ and $t_{\delta; v}$ is the $1 - \delta$ percentile t -value with v degrees of freedom. This method is referred to as EXCW method.

3.2 Confidence interval on γ_i using between group OLS estimator

From the fact that $Z'Z/l_2 = I_{l_1}$ and (2.4), \overline{Y} are normally distributed

$$\overline{Y} \sim N(X_1 \beta + X_2^* \gamma, (\sigma_\delta^2 + (\sigma_\epsilon^2/l_2))I_{l_1}) \tag{3.7}$$

The OLS estimators for β and γ from the between model are denoted by $\hat{\beta}_A$ and $\hat{\gamma}_A$, respectively, and they are obtained from (2.4). The estimators $\hat{\beta}_A$ and $\hat{\gamma}_A$ are called the between group OLS estimators for the first and second stage predictor variables, respectively, since they ignore variation within group. Normal equations for (2.4) are

$$\begin{pmatrix} X_1' \\ X_2^{*'} \end{pmatrix} (X_1 \ X_2^*) \begin{pmatrix} \hat{\beta}_A \\ \hat{\gamma}_A \end{pmatrix} = \begin{pmatrix} X_1' \\ X_2^{*'} \end{pmatrix} \overline{Y}. \tag{3.8}$$

From normal equations (3.8) and Formula (78) (Searle 1987, p263) $\hat{\gamma}_A$ are expressed as p_2 elements of the vector

$$\hat{\gamma}_A = (X_2^{*'} M_1 X_2^*)^{-1} X_2^{*'} M_1 \overline{Y} \tag{3.9}$$

where $M_1 = (I_{l_1} - X_1 (X_1' X_1)^{-1} X_1')$. By elementary algebra $\hat{\gamma}_A$ are normally distributed

$$\hat{\gamma}_A \sim N(\gamma, (\sigma_\delta^2 + (\sigma_\epsilon^2/l_2)) [X_2^* M_1 X_2^*]^{-1}). \tag{3.10}$$

Theorem 3.2 Under the distributional assumptions in (2.1), $\hat{\gamma}_A$ and R_1 are independent.

Proof R_1 is written as $R_1 = \underline{Y}'(Z/l_2)(I_{l_1} - H_1)(Z'/l_2)\underline{Y} = \overline{Y}'(I_{l_1} - H_1)\overline{Y}$. Note that $X_2^* M_1 (\sigma_\delta^2 + (\sigma_\epsilon^2/l_2)) I_{l_1} (I_{l_1} - H_1) = 0$ since $X_1' H_1 = X_1'$ and $X_2^* H_1 = X_2^*$. Therefore, $\hat{\gamma}_A$ and R_1 are independent.

Note that $(I_{l_1} - H_1)$ in R_1 is an idempotent matrix. Similar reasoning in (3.5) is applied to obtain an exact confidence interval. From (3.10) and Theorem 3.2 the between group OLS estimators of regression coefficients associated with secondary units in model are t distributed. Therefore, the proposed exact $1 - 2\alpha$ two-sided confidence interval on γ_i is

$$\hat{\gamma}_{iA} \pm t_{\alpha; l_1 - p_1 - p_2 - 1} \sqrt{\overline{X_2^*}^{\ddot{u}} S_\delta^2} \tag{3.11}$$

where $\hat{\gamma}_{iA}$ is the i th OLS estimator of regression coefficients in the secondary unit based on between group model and $\overline{X_2^*}^{\ddot{u}}$ is the i th diagonal element of matrix $(X_2^* M_1 X_2^*)^{-1}$. This method is referred to as EXCA method.

3.3 Confidence interval on β_i using between group OLS estimator

From normal equations (3.8) and formula (77) (Searle 1987, p263) $\hat{\beta}_A$ are expressed as $(1 + p_1)$ elements of the vector

$$\hat{\beta}_A = (X_1' X_1)^{-1} X_1' (\overline{Y} - X_2^* \hat{\gamma}_A). \tag{3.12}$$

From (3.7) and (3.10) $Var(\hat{\beta}_A)$ is written as

$$\begin{aligned} Var(\hat{\beta}_A) &= (X_1' X_1)^{-1} X_1' \{Var(\overline{Y}) + Var(X_2^* \hat{\gamma}_A) - 2COV(\overline{Y}, X_2^* \hat{\gamma}_A)\} X_1 (X_1' X_1)^{-1} \\ &= (X_1' X_1)^{-1} X_1' \{Var(\overline{Y}) + Var(X_2^* \hat{\gamma}_A)\} X_1 (X_1' X_1)^{-1} \\ &\quad - 2(X_1' X_1)^{-1} X_1' COV(\overline{Y}, X_2^* (X_2^* M_1 X_2^*)^{-1} X_2^* M_1 \overline{Y}) X_1 (X_1' X_1)^{-1} \\ &= (X_1' X_1)^{-1} X_1' \{(\sigma_\delta^2 + (\sigma_\epsilon^2/l_2)) \\ &\quad + (\sigma_\delta^2 + (\sigma_\epsilon^2/l_2)) X_2^* (X_2^* M_1 X_2^*)^{-1} X_2^*\} X_1 (X_1' X_1)^{-1} \\ &\quad - 2(\sigma_\delta^2 + (\sigma_\epsilon^2/l_2)) (X_1' X_1)^{-1} X_1' M_1 X_2^* (X_2^* M_1 X_2^*)^{-1} X_2^* X_1 (X_1' X_1)^{-1} \\ &= (\sigma_\delta^2 + (\sigma_\epsilon^2/l_2)) (X_1' X_1)^{-1} \\ &\quad + (X_1' X_1)^{-1} X_1' X_2^* [X_2^* M_1 X_2^*]^{-1} X_2^* X_1 (X_1' X_1)^{-1} \end{aligned}$$

by using that M_1 is symmetric and $X_1' M_1 = 0$. Since $\hat{\beta}_A$ are also a linear combination of \underline{Y} , $\hat{\beta}_A$ are normally distributed

$$\hat{\beta}_A \sim N(\beta, V(\hat{\beta}_A)) \quad (3.13)$$

where

$$V(\hat{\beta}_A) = (\sigma_\delta^2 + (\sigma_\epsilon^2/l_2))\{(X_1' X_1)^{-1} + (X_1' X_1)^{-1} X_1' X_2^* [X_2^{*'} M_1 X_2^*]^{-1} X_2^{*'} X_1 (X_1' X_1)^{-1}\}.$$

Theorem 3.3 Under the distributional assumptions in (2.1), $\hat{\beta}_A$ and R_1 are independent.

Proof $\hat{\beta}_A$ are written as $\hat{\beta}_A = (X_1' X_1)^{-1} X_1' (\bar{Y} - X_2^{*'} \hat{\gamma}_A) = (X_1' X_1)^{-1} X_1' (I_{l_1} - X_2^{*'} [X_2^{*'} M_1 X_2^*]^{-1} X_2^{*'} M_1) \bar{Y}$ by using (3.9). Note that $(X_1' X_1)^{-1} X_1' (I_{l_1} - X_2^{*'} [X_2^{*'} M_1 X_2^*]^{-1} X_2^{*'} M_1) (\sigma_\delta^2 + (\sigma_\epsilon^2/l_2)) I_{l_1} (I_{l_1} - H_1) = 0$ since $X_1' H_1 = X_1'$, $(X_1' X_1)^{-1} X_1' (I_{l_1} - H_1) = 0$, $M_1 (I_{l_1} - H_1) = I_{l_1} - H_1$, and $X_2^{*'} H_1 = X_2^{*}'$. Thus $\hat{\beta}_A$ and R_1 are independent.

Theorem 3.3 and (3.13) are used to construct a confidence interval on the first stage regression coefficients in the model. Similar reasoning in (3.5) is applied to obtain an exact confidence interval. The between group OLS estimators of regression coefficients associated with primary units in the model are t distributed. Thus, the proposed exact $1-2\alpha$ two-sided confidence interval on β_i is

$$\hat{\beta}_{iA} \pm t_{\alpha; l_1 - p_1 - p_2 - 1} \sqrt{X^{* \prime i} S_\delta^2} \quad (3.14)$$

where $\hat{\beta}_{iA}$ is the i th OLS estimator of regression coefficients in the primary units based on between group model, $X^{* \prime i}$ is the i th diagonal element of matrix $\{(X_1' X_1)^{-1} + (X_1' X_1)^{-1} X_1' X_2^{*'} [X_2^{*'} M_1 X_2^*]^{-1} X_2^{*'} X_1 (X_1' X_1)^{-1}\}$, and $S_\delta^2 = R_1 / (l_1 - p_1 - p_2 - 1)$. This method is referred to as EXB method.

3.4 Confidence interval on γ , using within group GLS estimator

The GLS estimators for $\underline{\gamma}$ from the within model, $\hat{\gamma}_{WG}$, are expressed as p_2 elements of the vector and are written as

$$\hat{\gamma}_{WG} = (\bar{X}_2' W^{-1} \bar{X}_2)^{-1} \bar{X}_2' W^{-1} \bar{Y} \quad (3.15)$$

since $V_{\bar{Y}} = \text{Var}(\bar{Y}) = \sigma_\epsilon^2 W$. In addition, the GLS estimators, $\hat{\gamma}_{WG}$, are normally distributed

$$\hat{\gamma}_{WG} \sim N(\gamma, \sigma_\epsilon^2 (\bar{X}_2' W^{-1} \bar{X}_2)^{-1}). \quad (3.16)$$

Theorem 3.4 Under the distributional assumptions in (2.1), $\hat{\gamma}_{WG}$ and R_2 are independent.

Proof Note that $R_2 = \underline{Y}'(I_{l_2} - H_2)\underline{Y}$. Thus, $(\overline{X}_2' W^{-1} \overline{X}_2)^{-1} \overline{X}_2' W^{-1} (\sigma_\epsilon^2 W)$
 $(I_{l_2} - H_2) = 0$ since $\overline{X}_2' H_2 = \overline{X}_2'$. It follows that $\hat{\gamma}_{WG}$ and R_2 are independent.

Similar reasoning in (3.5) is applied to obtain an exact confidence interval. The GLS estimators of regression coefficients associated with secondary units in the model are t distributed from (3.16) and Theorem 3.4. Therefore, the proposed exact $1 - 2\alpha$ two-sided confidence interval on γ_i is

$$\hat{\gamma}_{iWG} \pm t_{\alpha; l_1 l_2 - l_1 - l_2} \sqrt{\overline{X}_w^{*ii} S_\epsilon^2} \tag{3.17}$$

where $\hat{\gamma}_{iWG}$ is the i th GLS estimator of regression coefficients in the secondary unit based on within model and \overline{X}_w^{*ii} is the i th diagonal element of matrix $(\overline{X}_2' W^{-1} \overline{X}_2)^{-1}$. This method is referred to as EXCG method.

The GLS estimators for $\underline{\beta}$ and $\underline{\gamma}$ from the between model are same as $\hat{\beta}_A$ and $\hat{\gamma}_A$, respectively, since $Var(\overline{Y}) = (\sigma_\delta^2 + (\sigma_\epsilon^2/l_2))I_{l_1}$ in (3.7). In particular, The GLS estimators for $\underline{\gamma}$ from the between model, $\hat{\gamma}_{AG}$, are written

$$\hat{\gamma}_{AG} = (X_2^{*'} M_1 X_2^*)^{-1} X_2^{*'} M_1 \overline{Y} \tag{3.18}$$

which are same as $\hat{\gamma}_A$ in (3.9). Similarly, the GLS estimators for $\underline{\beta}$ from the between model, $\hat{\beta}_{AG}$, are equivalent to $\hat{\beta}_A$ in (3.12).

4. Conclusions

This paper used OLS and GLS estimators to find regression coefficients in the model. The estimators, $\hat{\beta}$ and $\hat{\gamma}$, turned out to be independent of quadratic forms, R_1 and R_2 , obtained from the model. One exact two-sided confidence interval on β , and three exact two-sided confidence interval on γ_i were proposed. EXB method is recommended to construct confidence interval on regression coefficients associated with the primary unit in the model. EXCW, EXCA, and EXCG methods are recommended to construct confidence intervals on regression coefficients associated with the secondary unit in the model. The proposed confidence intervals on regression coefficients could be used in two stage regression model application areas.

References

- [1] Christensen, R. (1987), The Analysis of Two-Stage Sampling Data by Ordinary Least Squares, *Journal of American Statistical Association*, 82, 492-498.
- [2] Fuller, W. A. and Battese, G. E. (1973), Transformations for Estimation of Linear Models with Nested-Error Structure, *Journal of American Statistical Association*, 68, 626-632.
- [3] Park, D. J. and Burdick, R. K. (1994), Confidence intervals on the regression coefficients in a simple linear regression model with a balanced one-fold nested error structure, *Communications in Statistics- Simulation and Computation*, 23(1), 43-58.
- [4] Park, D. J. (1996), Confidence intervals on the variance components in two stage regression model, *Journal of Statistical Theory and Methods*, 7, 87-92.
- [5] Searle, S. R. (1971), *Linear Models*, New York: John Wiley & Sons Inc.
- [6] Searle, S. R. (1987), *Linear Models For Unbalanced Data*, New York: John Wiley & Sons Inc.
- [7] Weerakkody, G. J. and Johnson, D. E. (1992), Estimation of Within Model Parameters in Regression Models With a Nested Error Structure, *Journal of American Statistical Association*, 87, 708-713.